

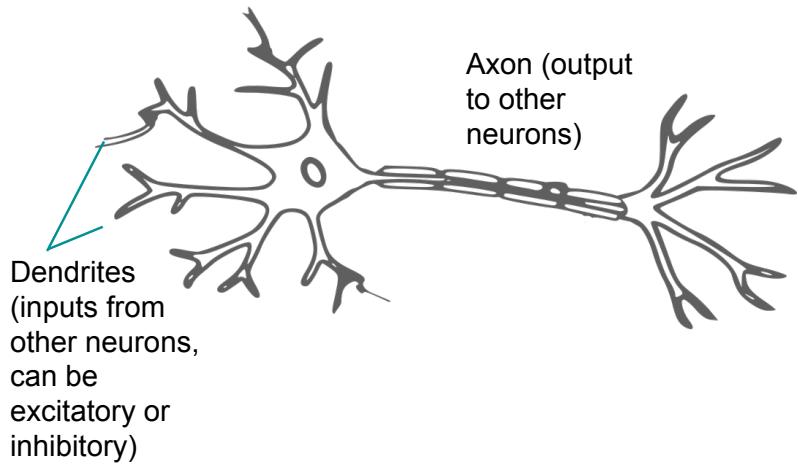
10-701 Introduction to Machine Learning (PhD) Lecture 9: Neural Networks

Leila Wehbe
Carnegie Mellon University
Machine Learning Department

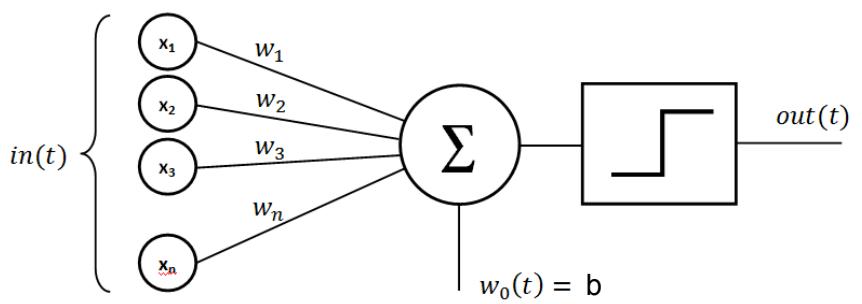
Slides based on Tom Mitchell's
10-701 Spring 2016 material
Readings: Hal Daumé III Chapter 4, 10
[TM] Ch. 4, [CB] Ch. 5

The Perceptron

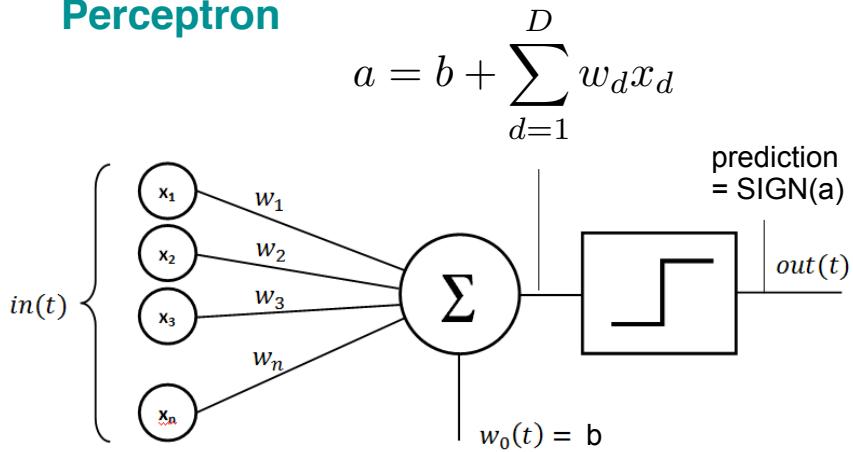
Inspired by biological neurons



Perceptron



Perceptron



Error driven learning

$$a = b + \sum_{d=1}^D w_d x_d$$

- At each step, return $\text{SIGN}(a)$
- if $\text{SIGN}(a) \neq y$ update parameter
- otherwise don't change

Algorithm 5 PERCEPTRONTRAIN($D, MaxIter$)

```

1:  $w_d \leftarrow 0$ , for all  $d = 1 \dots D$                                 // initialize weights
2:  $b \leftarrow 0$                                                         // initialize bias
3: for iter = 1 ... MaxIter do
4:   for all  $(x,y) \in D$  do
5:      $a \leftarrow \sum_{d=1}^D w_d x_d + b$                                 // compute activation for this example
6:     if  $y_a \leq 0$  then
7:        $w_d \leftarrow w_d + yx_d$ , for all  $d = 1 \dots D$                 // update weights
8:        $b \leftarrow b + y$                                               // update bias
9:     end if
10:   end for
11: end for
12: return  $w_0, w_1, \dots, w_D, b$ 

```

from http://ciml.info/dl/v0_99/ciml-v0_99-ch08.pdf

Does this move a in the right direction?

- update $\mathbf{w}' = \mathbf{w} + y\mathbf{x} = \mathbf{w} + \mathbf{x}$
- $b' = b + y = b+1$

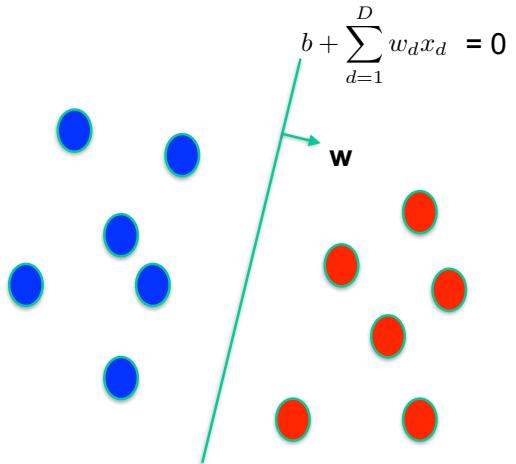
$$\begin{aligned}
 a' &= \sum_{d=1}^D w'_d x_d + b' \\
 &= \sum_{d=1}^D (w_d + x_d)x_d + (b + 1) \\
 &= \sum_{d=1}^D w_d x_d + b + \sum_{d=1}^D x_d x_d + 1 \\
 &= a + \sum_{d=1}^D x_d^2 + 1 \quad > \quad a
 \end{aligned}$$

Does this move \mathbf{w} in the right direction?

- update $\mathbf{w}' = \mathbf{w} + y\mathbf{x} = \mathbf{w} + \mathbf{x}$
- $b' = b + y = b+1$

$$\begin{aligned}a' &= \sum_{d=1}^D w'_d x_d + b' \\&= \sum_{d=1}^D (w_d + x_d) x_d + (b + 1) \\&= \sum_{d=1}^D w_d x_d + b + \sum_{d=1}^D x_d x_d + 1 \\&= a + \sum_{d=1}^D x_d^2 + 1 > a \quad \text{a becomes more positive (not guaranteed that } a>0)\end{aligned}$$

What is the decision boundary?



How good is this algorithm?

- Convergence: an entire pass without changing the weights.
- If the data is linearly separable, the algorithm will converge. But not necessarily to the “best” boundary

Notion of margin

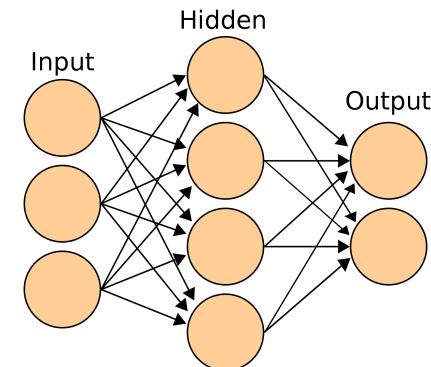
$$\text{margin}(\mathbf{D}, \mathbf{w}, b) = \begin{cases} \min_{(x,y) \in \mathbf{D}} y(\mathbf{w} \cdot \mathbf{x} + b) & \text{if } \mathbf{w} \text{ separates } \mathbf{D} \\ -\infty & \text{otherwise} \end{cases}$$

$$\text{margin}(\mathbf{D}) = \sup_{\mathbf{w}, b} \text{margin}(\mathbf{D}, \mathbf{w}, b)$$

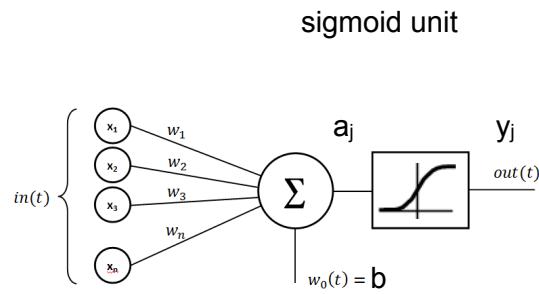
If data is linearly separable with margin γ and $\|\mathbf{x}\| \leq 1$, then algorithm will converge in $\frac{1}{\gamma^2}$ updates

Neural Networks

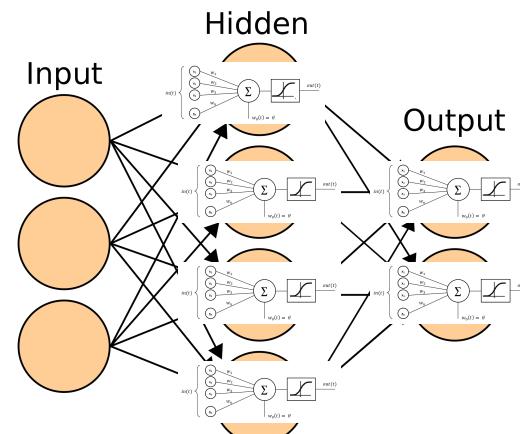
Every node is analogous to a neuron



Every node is analogous to a neuron

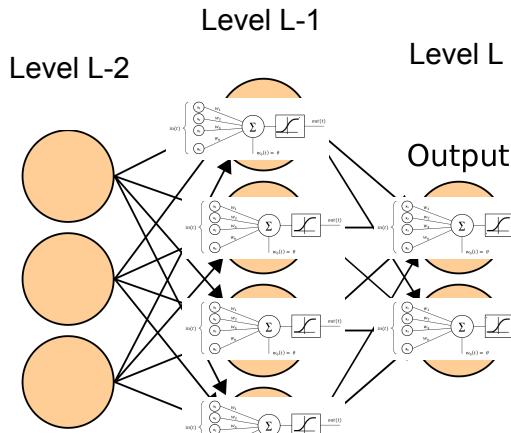


Every node is analogous to a neuron



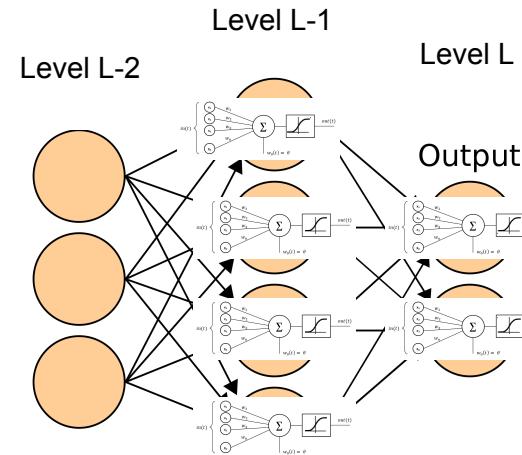
**Prediction:
Forward Propagation**

$$y_{L-1}^j = \sigma \left(\sum_k w_{L-1}^{jk} y_{L-2}^k + b_{L-1}^j \right)$$

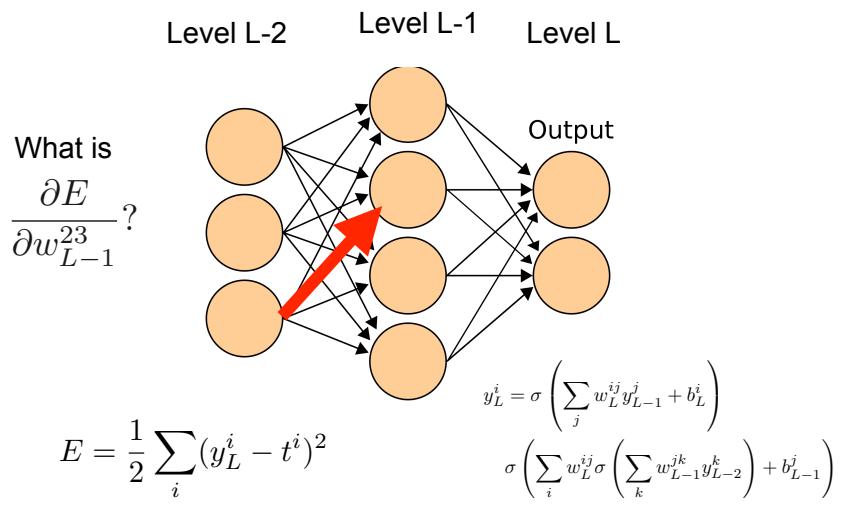


**Prediction:
Forward Propagation**

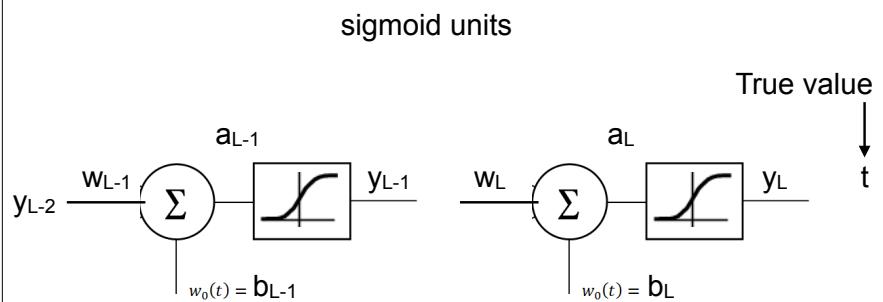
$$\begin{aligned} y_L^i &= \sigma \left(\sum_j w_L^{ij} y_{L-1}^j + b_L^i \right) \\ &\sigma \left(\sum_j w_L^{ij} \sigma \left(\sum_k w_{L-1}^{jk} y_{L-2}^k + b_{L-1}^j \right) + b_L^i \right) \end{aligned}$$



How to train?



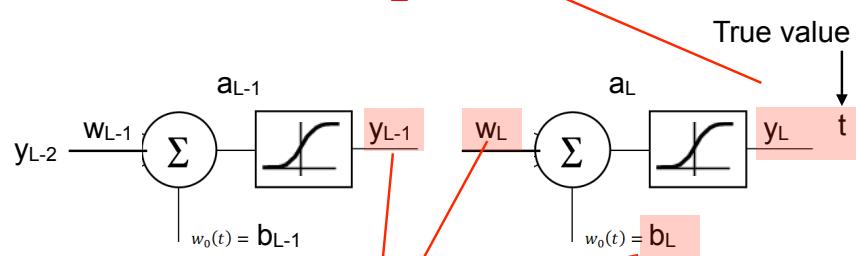
Backprop with one node per layer



Backprop with one node per layer

Loss function

$$E = \frac{1}{2}(y_L - t)^2$$

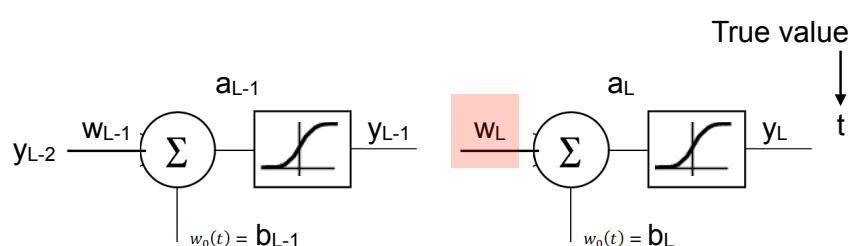


Find the gradient for one training point at earlier nodes and parameters: how much does a very small change in the value of nodes and parameters affect the loss for that point?

Backprop with one node per layer

Loss function

$$E = \frac{1}{2}(y_L - t)^2$$

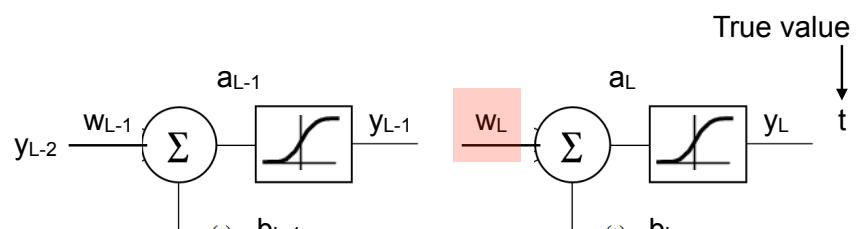


$$\frac{\partial E}{\partial w_L}$$

Backprop with one node per layer

Loss function

$$E = \frac{1}{2}(y_L - t)^2$$

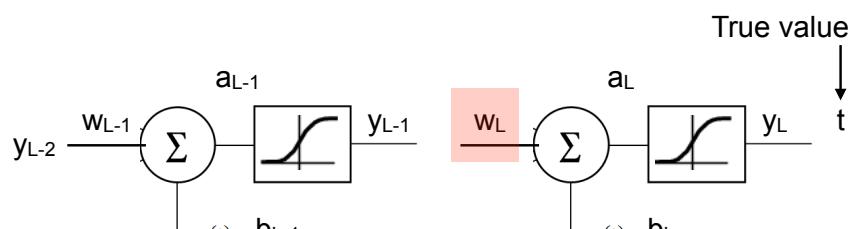


$$\frac{\partial E}{\partial w_L}$$

$$\begin{aligned} E &= \frac{1}{2}(y_L - t)^2 \\ y_L &= \sigma(a_L) \\ a_L &= w_L y_{L-1} + b_L \end{aligned}$$

Backprop with one node per layer

$$E = \frac{1}{2}(y_L - t)^2$$

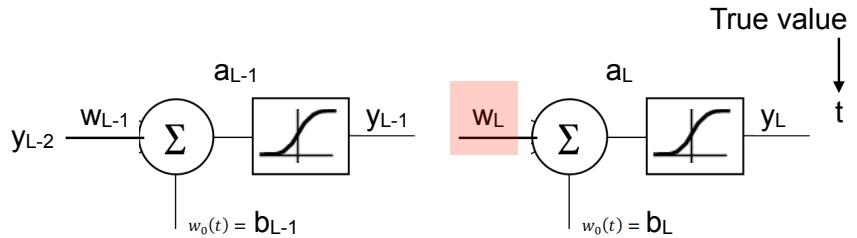


$$\frac{\partial E}{\partial w_L} = \frac{\partial a_L}{\partial w_L} \frac{\partial y_L}{\partial a_L} \frac{\partial E}{\partial y_L}$$

$$\begin{aligned} E &= \frac{1}{2}(y_L - t)^2 \\ y_L &= \sigma(a_L) \\ a_L &= w_L y_{L-1} + b_L \end{aligned}$$

Backprop with one node per layer

$$E = \frac{1}{2}(y_L - t)^2$$



$$\frac{\partial E}{\partial w_L} = \frac{\partial E}{\partial w_L} \frac{\partial y_L}{\partial a_L} \frac{\partial a_L}{\partial y_{L-1}}$$

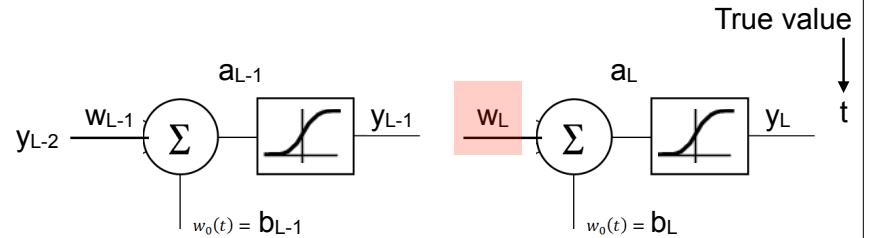
$E = \frac{1}{2}(y_L - t)^2$

$y_L = \sigma(a_L)$

$a_L = w_L y_{L-1} + b_L$

Backprop with one node per layer

$$E = \frac{1}{2}(y_L - t)^2$$



$$\frac{\partial E}{\partial w_L} = y_{L-1} \sigma'(a_L) (y_L - t)$$

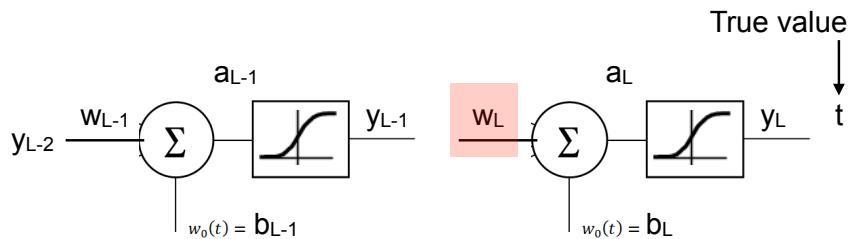
$E = \frac{1}{2}(y_L - t)^2$

$y_L = \sigma(a_L)$

$a_L = w_L y_{L-1} + b_L$

Backprop with one node per layer

$$E = \frac{1}{2}(y_L - t)^2$$



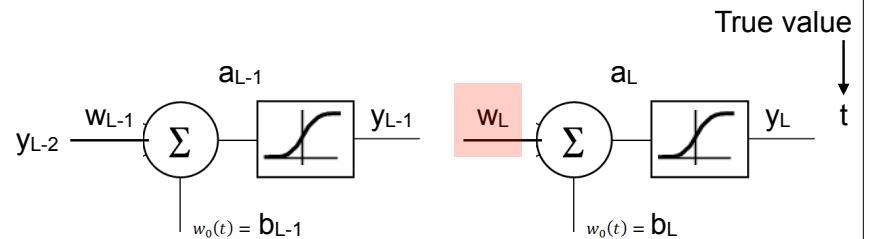
$$\frac{\partial E}{\partial w_L} = y_{L-1} \sigma'(a_L) (y_L - t)$$

$\sigma'(a_L) = \sigma(a_L)(1 - \sigma(a_L))$

$\sigma(a_L) = y_L$

Backprop with one node per layer

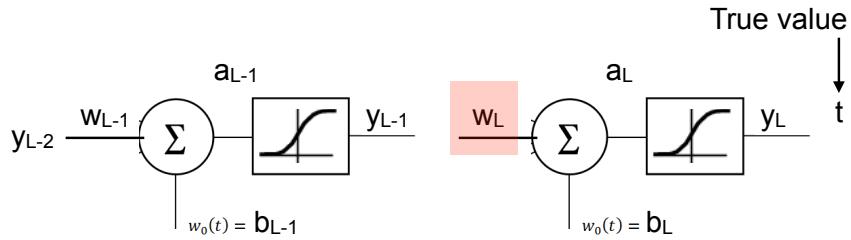
$$E = \frac{1}{2}(y_L - t)^2$$



$$\frac{\partial E}{\partial w_L} = y_{L-1} y_L (1 - y_L) (y_L - t)$$

Backprop with one node per layer

$$E = \frac{1}{2}(y_L - t)^2$$

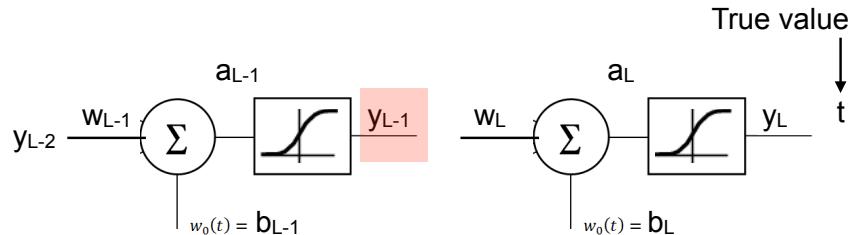


Effect of ∂w_L depends on y_{L-1}

$$\frac{\partial E}{\partial w_L} = y_{L-1} y_L (1 - y_L) (y_L - t)$$

Backprop with one node per layer

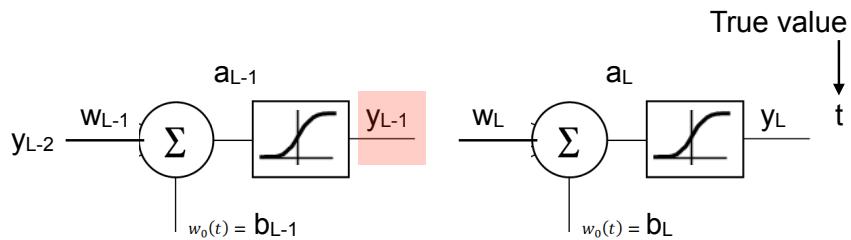
$$E = \frac{1}{2}(y_L - t)^2$$



$$\frac{\partial E}{\partial y_{L-1}} = \frac{\partial a_L}{\partial y_{L-1}} \frac{\partial y_L}{\partial a_L} \frac{\partial E}{\partial y_L}$$

Backprop with one node per layer

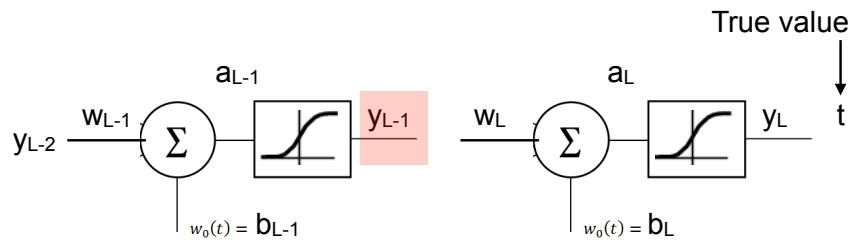
$$E = \frac{1}{2}(y_L - t)^2$$



$$\frac{\partial E}{\partial y_{L-1}} = \frac{\partial a_L}{\partial y_{L-1}} \frac{\partial y_L}{\partial a_L} \frac{\partial E}{\partial y_L} \quad a_L = w_L y_{L-1} + b_L$$

Backprop with one node per layer

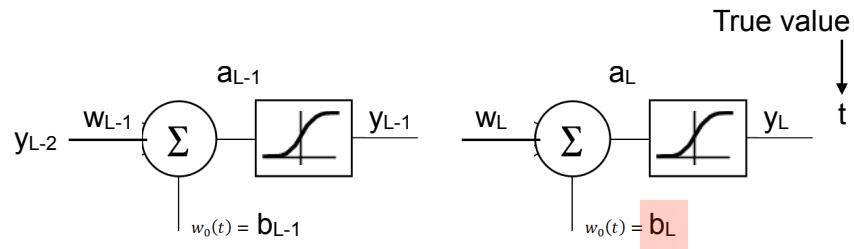
$$E = \frac{1}{2}(y_L - t)^2$$



$$\frac{\partial E}{\partial y_{L-1}} = w_L y_L (1 - y_L) (y_L - t) \quad a_L = w_L y_{L-1} + b_L$$

Backprop with one node per layer

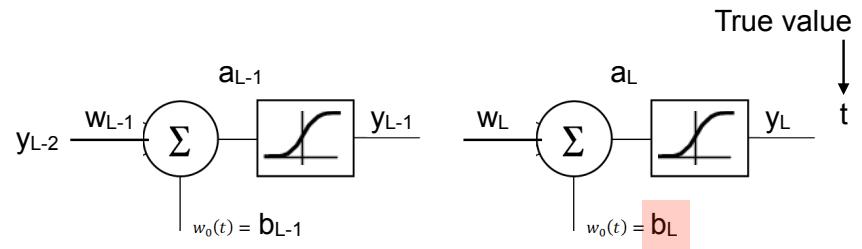
$$E = \frac{1}{2}(y_L - t)^2$$



$$\frac{\partial E}{\partial b_L} = \frac{\partial a_L}{\partial b_L} \frac{\partial y_L}{\partial a_L} \frac{\partial E}{\partial y_L}$$

Backprop with one node per layer

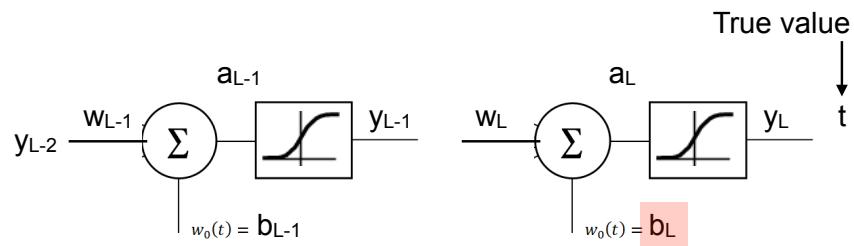
$$E = \frac{1}{2}(y_L - t)^2$$



$$\frac{\partial E}{\partial b_L} = \frac{\partial a_L}{\partial b_L} \frac{\partial y_L}{\partial a_L} \frac{\partial E}{\partial y_L} \quad a_L = w_L y_{L-1} + b_L$$

Backprop with one node per layer

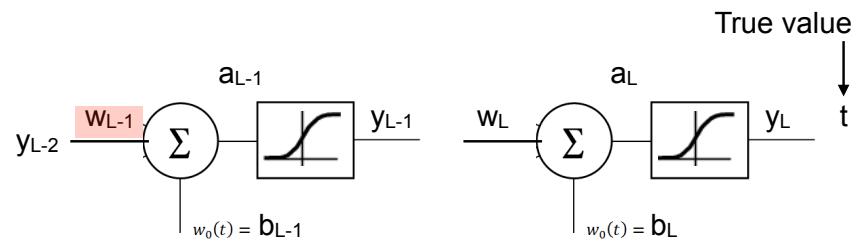
$$E = \frac{1}{2}(y_L - t)^2$$



$$\frac{\partial E}{\partial b_L} = \frac{1}{2} y_L (1 - y_L) (y_L - t)$$

Backprop with one node per layer

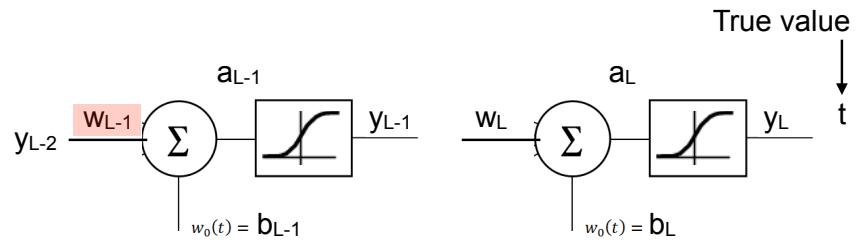
$$\text{Loss function } E = \frac{1}{2}(y_L - t)^2$$



$$\frac{\partial E}{\partial w_{L-1}}$$

Backprop with one node per layer

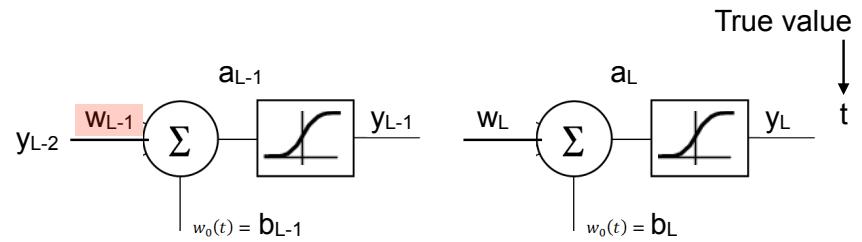
Loss function $E = \frac{1}{2}(y_L - t)^2$



$$\frac{\partial E}{\partial w_{L-1}} = \frac{\partial a_{L-1}}{\partial w_{L-1}} \frac{\partial y_{L-1}}{\partial a_{L-1}} \frac{\partial E}{\partial y_{L-1}}$$

Backprop with one node per layer

Loss function $E = \frac{1}{2}(y_L - t)^2$



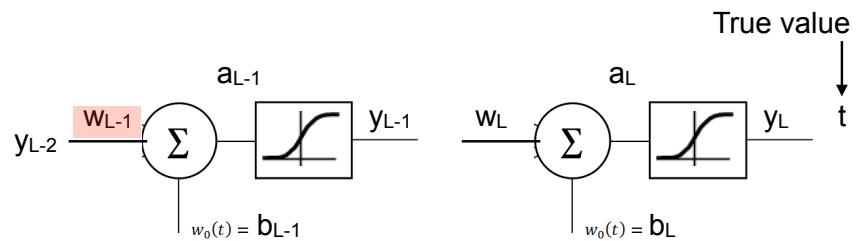
$$a_{L-1} = w_{L-1}y_{L-2} + b_{L-1} \quad y_{L-1} = \sigma(a_{L-1})$$

$$\frac{\partial E}{\partial w_{L-1}} = \frac{\partial a_{L-1}}{\partial w_{L-1}} \frac{\partial y_{L-1}}{\partial a_{L-1}} \frac{\partial E}{\partial y_{L-1}}$$

Already Computed

Backprop with one node per layer

Loss function $E = \frac{1}{2}(y_L - t)^2$

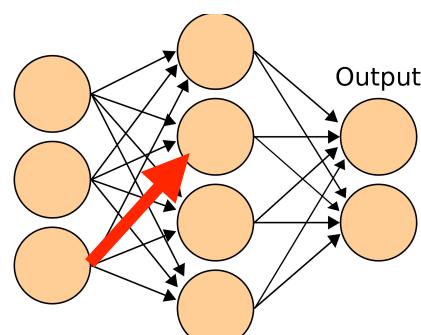


$$\frac{\partial E}{\partial w_{L-1}} = y_{L-2}y_{L-1}(1 - y_{L-1})w_Ly_L(1 - y_L)(y_L - t)$$

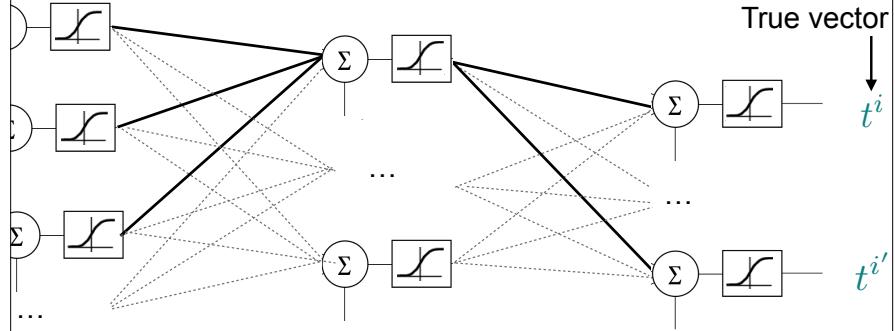
Backprop with more than one node per layer

Level L-2 Level L-1 Level L

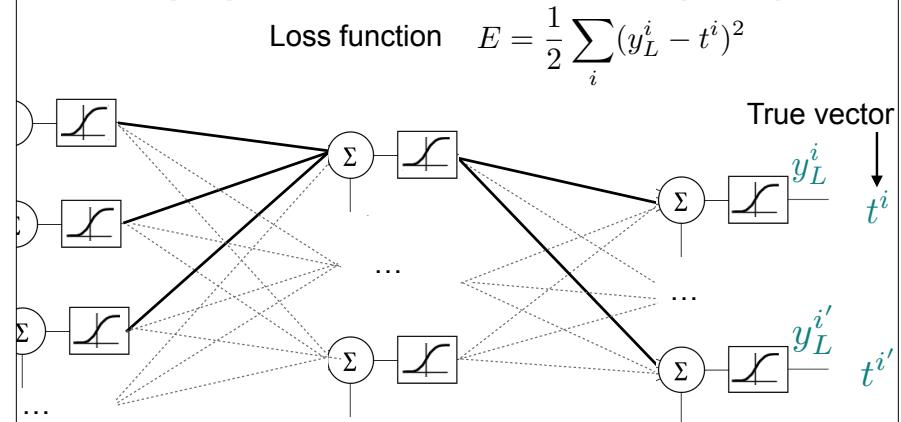
What is
 $\frac{\partial E}{\partial w_{L-1}^{23}}$?



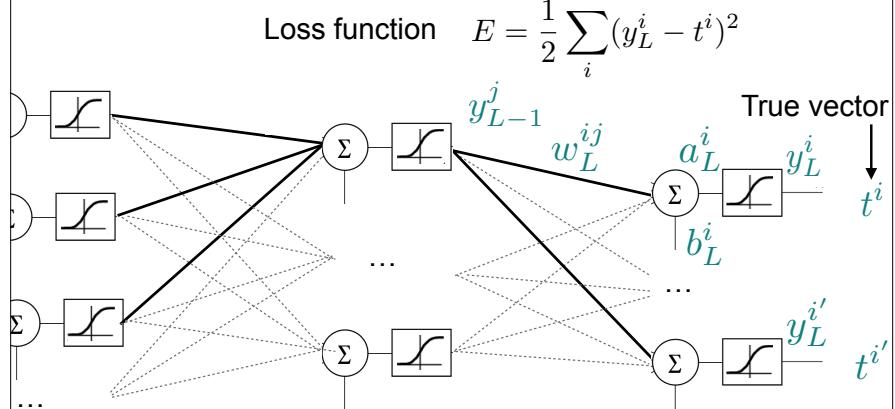
Backprop with more than one node per layer



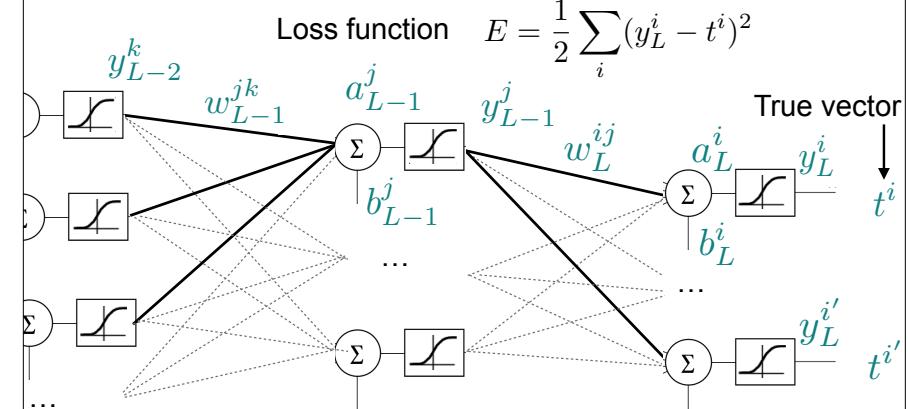
Backprop with more than one node per layer



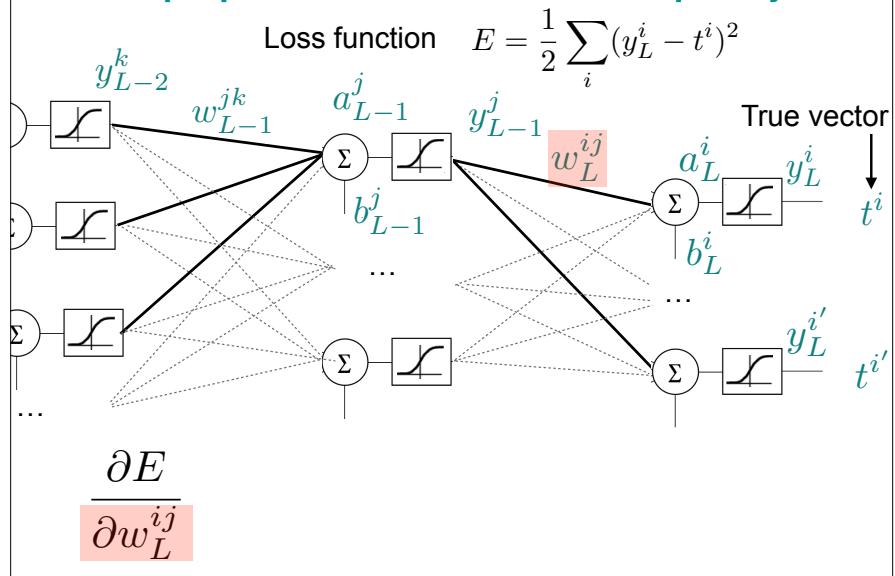
Backprop with more than one node per layer



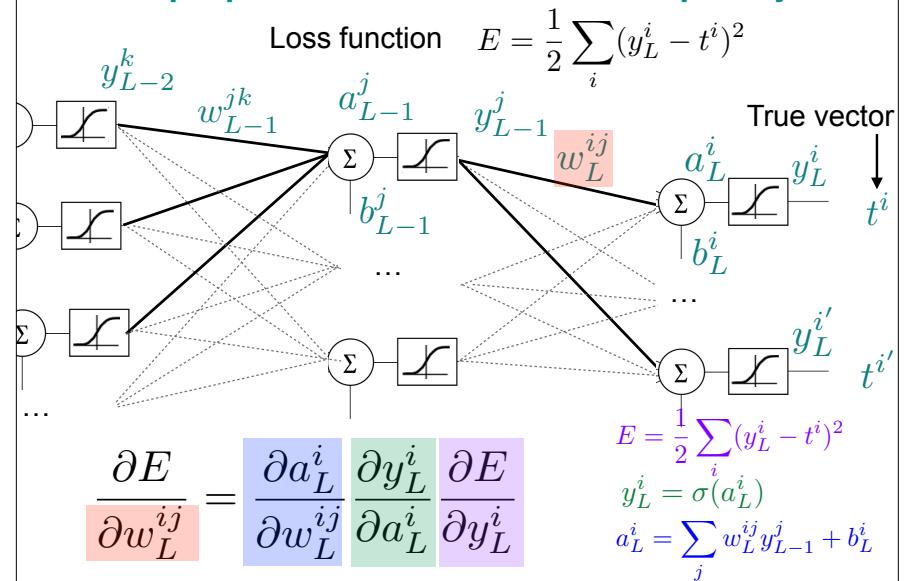
Backprop with more than one node per layer



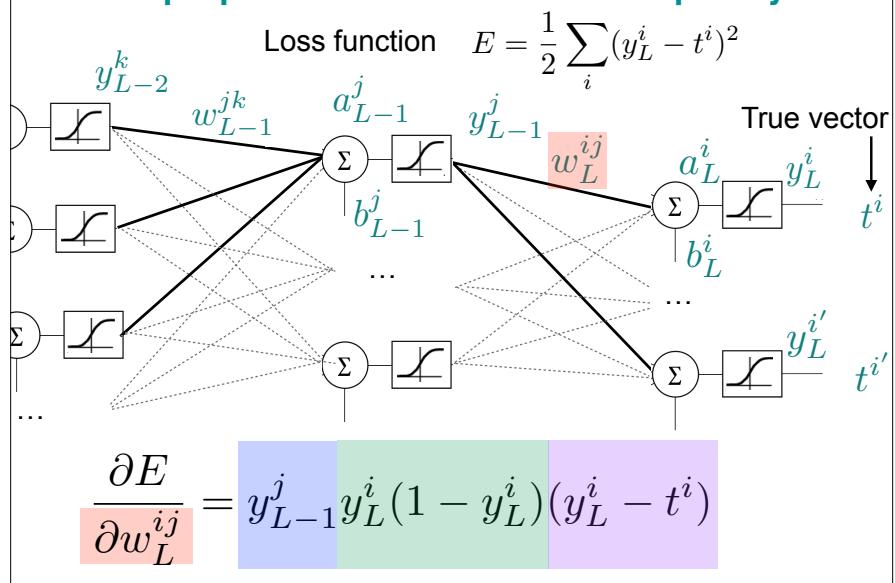
Backprop with more than one node per layer



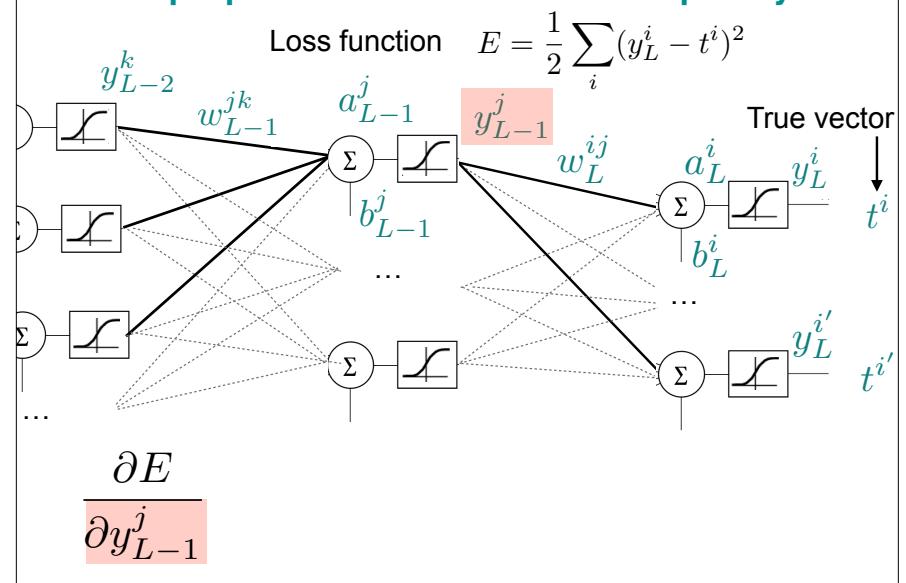
Backprop with more than one node per layer



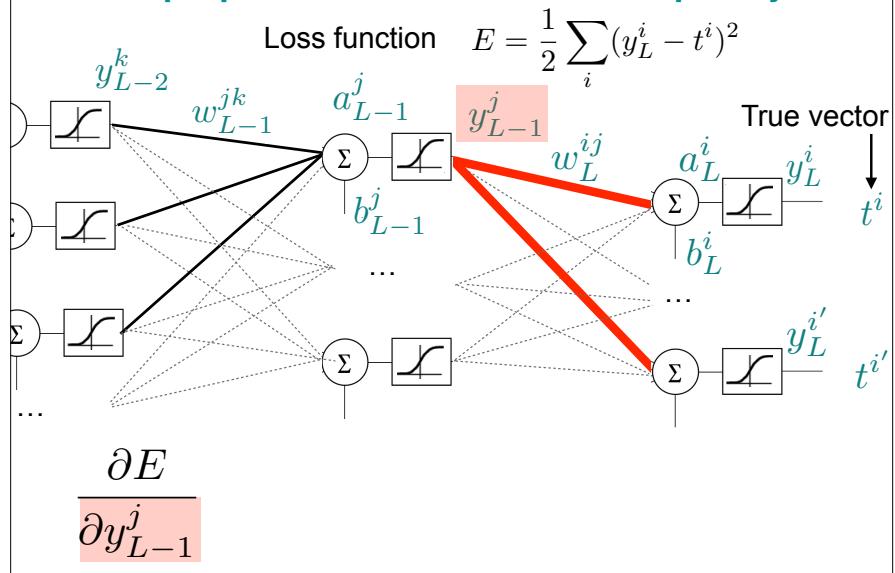
Backprop with more than one node per layer



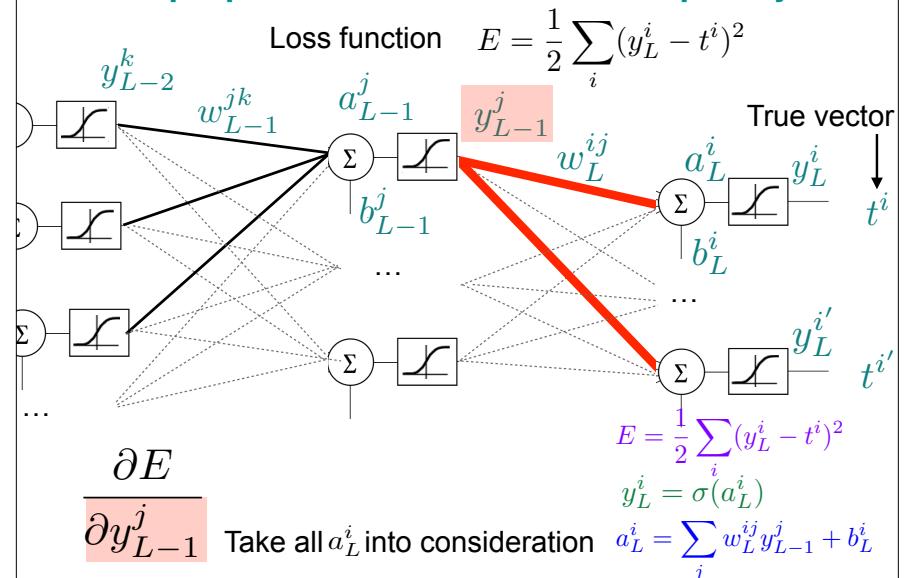
Backprop with more than one node per layer



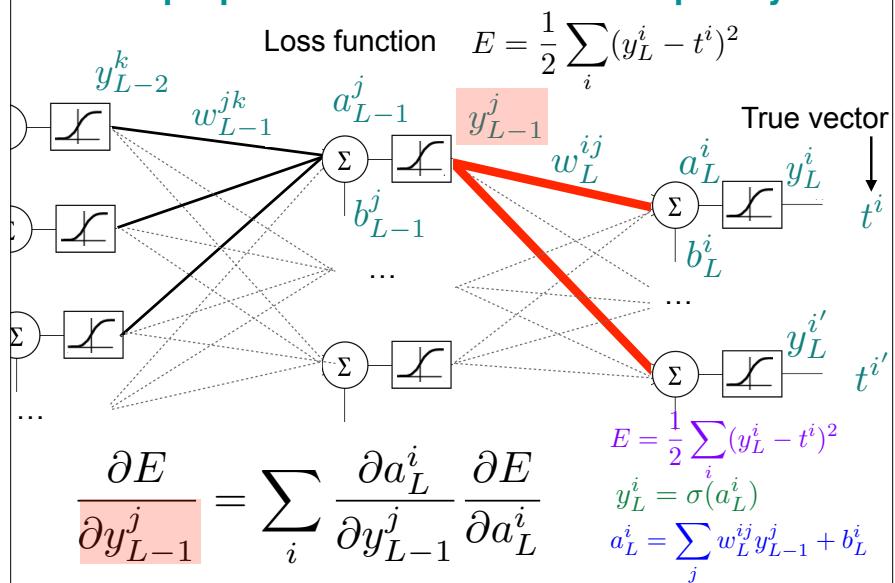
Backprop with more than one node per layer



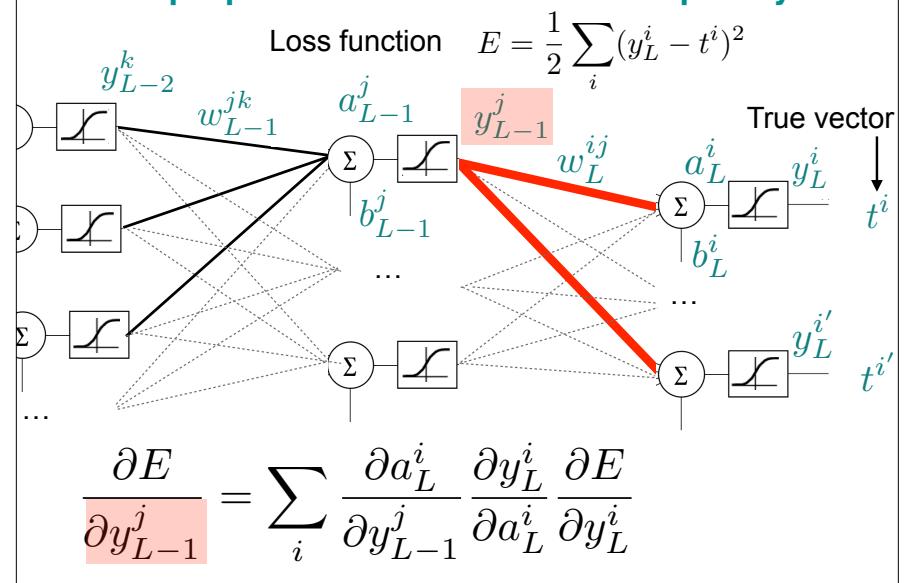
Backprop with more than one node per layer



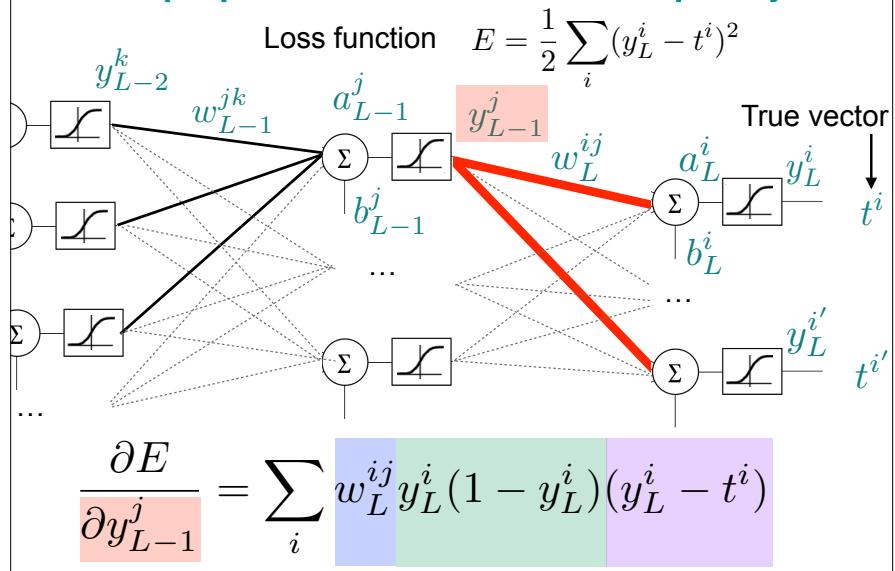
Backprop with more than one node per layer



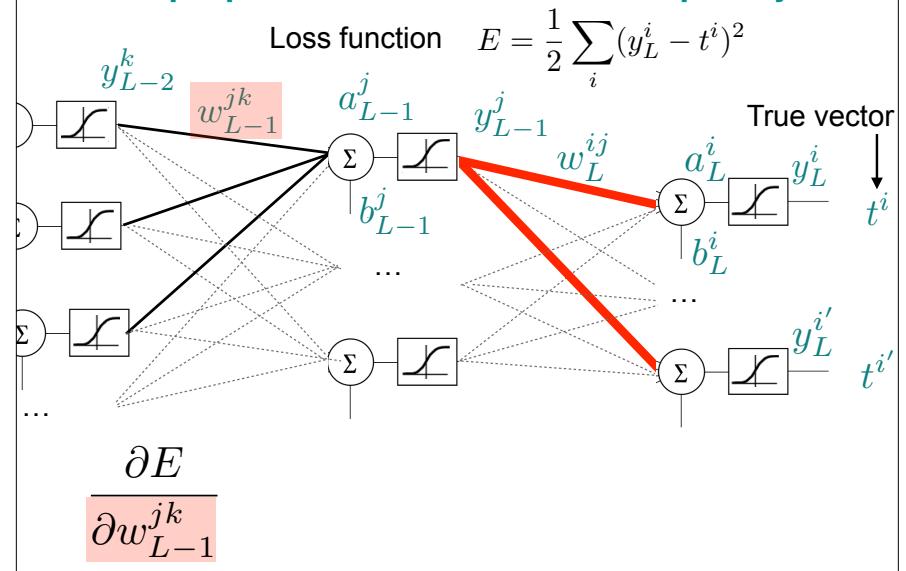
Backprop with more than one node per layer



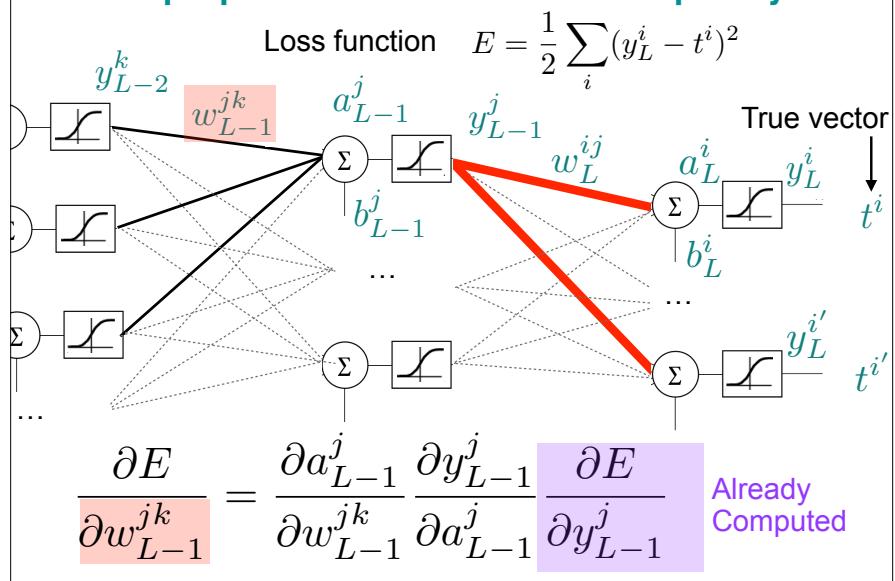
Backprop with more than one node per layer



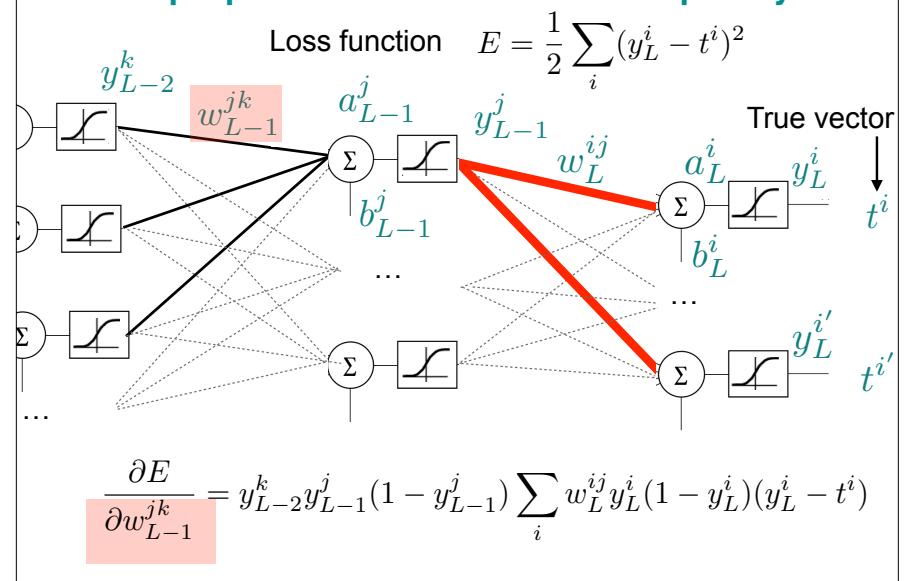
Backprop with more than one node per layer



Backprop with more than one node per layer



Backprop with more than one node per layer



How to train

Random initialization, then repeat:

- Sample a training example
 - Input training example and compute all outputs y from all nodes

How to train

Random initialization, then repeat:

- Sample a training example
 - Input training example and compute all outputs y from all nodes
 - Back-propagate E (weighting it by the gradient of previous layer and activation function), i.e. compute $\frac{\partial E}{\partial y_H^j}$ for all layers starting from output layer L to first layer 1.

How to train

Random initialization, then repeat:

- Sample a training example
 - Input training example and compute all outputs y from all nodes
 - Back-propagate E (weighting it by the gradient of previous layer and activation function), i.e. compute $\frac{\partial E}{\partial y_H^j}$ for all layers starting from output layer L to first layer 1.
- Calculate the gradients $\frac{\partial E}{\partial w_H^{ij}}$ and $\frac{\partial E}{\partial b_H^i}$

How to train

Random initialization, then repeat:

- Sample a training example
 - Input training example and compute all outputs y from all nodes
 - Back-propagate E (weighting it by the gradient of previous layer and activation function), i.e. compute $\frac{\partial E}{\partial y_H^j}$ for all layers starting from output layer L to first layer 1.
- Calculate the gradients $\frac{\partial E}{\partial w_H^{ij}}$ and $\frac{\partial E}{\partial b_H^i}$
- Update the parameters
$$w_H^{ij} \leftarrow w_H^{ij} - \eta \frac{\partial E}{\partial w_H^{ij}} \quad b_H^i \leftarrow b_H^i - \eta \frac{\partial E}{\partial b_H^i}$$

How to train

Random initialization, then repeat:

- Sample a training example
- Input training example and compute all outputs y from all nodes
- Back-propagate E (weighting it by the gradient of previous layer and activation function), i.e. compute $\frac{\partial E}{\partial y_H^i}$ for all layers starting from output layer L to first layer 1.
- Calculate the gradients $\frac{\partial E}{\partial w_H^{ij}}$ and $\frac{\partial E}{\partial b_H^i}$
- Update the parameters

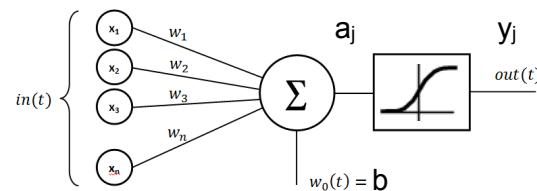
$$w_H^{ij} \leftarrow w_H^{ij} - \eta \frac{\partial E}{\partial w_H^{ij}} \quad b_H^i \leftarrow b_H^i - \eta \frac{\partial E}{\partial b_H^i}$$

but... function not convex...

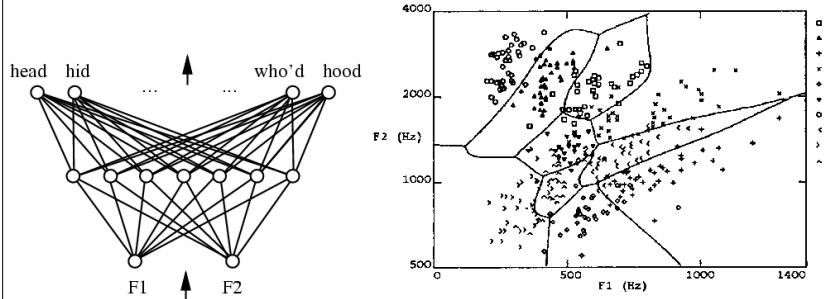
still gradient descent could lead to "good" solutions

Questions

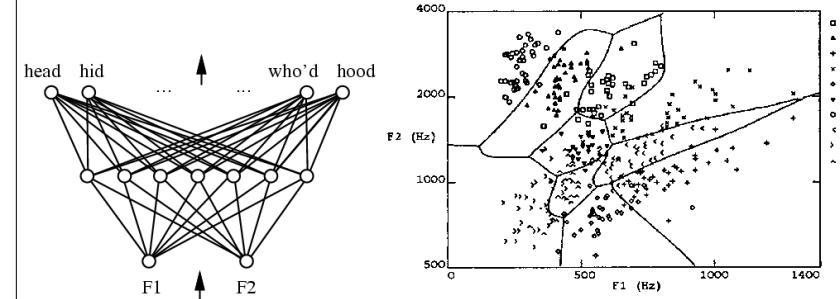
- Is there a difference between LR and one sigmoid unit?



Multilayer Networks of Sigmoid Units

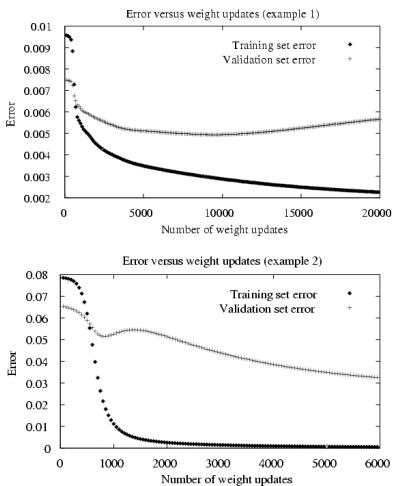


Multilayer Networks of Sigmoid Units



Can learn highly non-linear decision surfaces.
However, the function that a deep neural net learns is typically difficult to interpret

Overfitting in ANNs



Regularization

- To avoid overfitting, add penalty, e.g. L2 penalty:

$$\arg \min_W \sum_l \|\hat{f}(\mathbf{x}^l) - \mathbf{t}^l\|_2^2$$

$$\arg \min_W \sum_l \|\hat{f}(\mathbf{x}^l) - \mathbf{t}^l\|_2^2 + \lambda \|W\|_2^2$$

Regularization

- To avoid overfitting, add penalty, e.g. L2 penalty:

$$\arg \min_W \sum_l \|\hat{f}(\mathbf{x}^l) - \mathbf{t}^l\|_2^2$$

$$\arg \min_W \sum_l \|\hat{f}(\mathbf{x}^l) - \mathbf{t}^l\|_2^2 + \lambda \|W\|_2^2$$

- Early stopping: stop when validation performance worsens. This has a regularization effect

Regularization

- To avoid overfitting, add penalty, e.g. L2 penalty:

$$\arg \min_W \sum_l \|\hat{f}(\mathbf{x}^l) - \mathbf{t}^l\|_2^2$$

$$\arg \min_W \sum_l \|\hat{f}(\mathbf{x}^l) - \mathbf{t}^l\|_2^2 + \lambda \|W\|_2^2$$

- Early stopping: stop when validation performance worsens. This has a regularization effect
- Dropout: randomly ignore a proportion of nodes p during training (at each iteration). Account for this during testing (multiplying by p) or by scaling inversely during training (dividing by p).

How to train

- What happens if we initialize weights to 0?

How to train

- What happens if we initialize weights to 0?
- Random initialization:
 - Different results (we might learn the same function but at different nodes, or learn different functions since the function is not convex)
 - Repeat multiple times and pick the one with best cross-validation performance

Many hyperparameters

- Many knobs to tune:
 - Structure of network
 - Number of layers
 - Number of nodes per layer
 - Other: recurrence, convolution
 - Learning rate
 - Initialization
 - Loss function
 - Gate
 - Stopping criterion (early stopping?)

Many hyperparameters

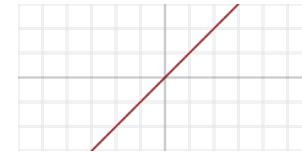
- Many knobs to tune:
 - Structure of network
 - Number of layers
 - Number of nodes per layer
 - Other: recurrence, convolution
 - Learning rate
 - Initialization
 - Loss function
 - Gate
 - Stopping criterion

Many additional tweaks that can help

- Batch normalization
- Learning rate decay
- Momentum
- ...

Choice of activation gate

Linear

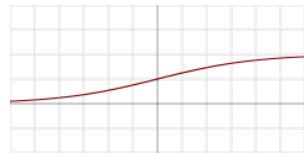


$$f(x) = x \qquad f'(x) = 1$$

Can be used to predict continuous values at output.
What happens when you stack linear layers?

Choice of activation gate

Sigmoid

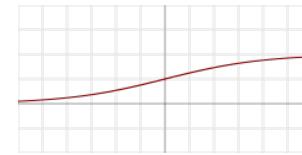


$$f(x) = \sigma(x) = \frac{1}{1 + e^{-x}} \qquad f'(x) = f(x)(1 - f(x))$$

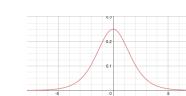
Outputs between 0 and 1, can be used for probability
Can saturate when very low or very high weights
Contributes to vanishing gradient

Choice of activation gate

Sigmoid



$f'(x)$

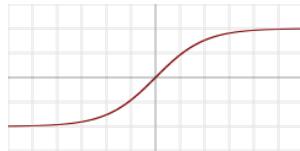


$$f(x) = \sigma(x) = \frac{1}{1 + e^{-x}} \qquad f'(x) = f(x)(1 - f(x))$$

Outputs between 0 and 1, can be used for probability
Can saturate
Contributes to vanishing gradient

Choice of activation gate

tanh

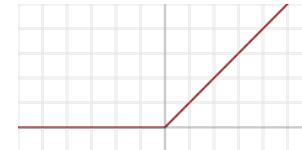


$$f(x) = \tanh(x) = \frac{(e^x - e^{-x})}{(e^x + e^{-x})} \quad f'(x) = 1 - f(x)^2$$

Range -1 to 1

Choice of activation gate

Rectified Linear Unit (ReLU)

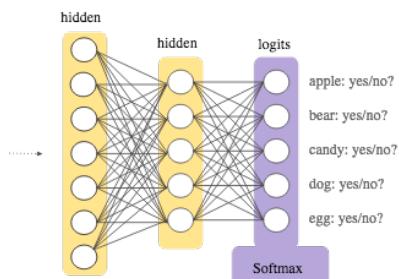


$$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases} \quad f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$$

Solve the vanishing gradient problem, but have a problem that nodes might die when negative value and never update. Can fix with leaky ReLU

Choice of activation gate

Softmax

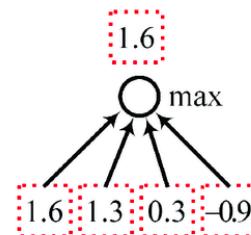


$$f_i(\vec{x}) = \frac{e^{x_i}}{\sum_{j=1}^J e^{x_j}}$$

$$\frac{\partial f_i(\vec{x})}{\partial x_j} = f_i(\vec{x})(\delta_{ij} - f_j(\vec{x}))$$

Choice of activation gate

Maxpool



$$f(\vec{x}) = \max_i x_i$$

$$\frac{\partial f}{\partial x_j} = \begin{cases} 1 & \text{for } j = \operatorname{argmax}_i x_i \\ 0 & \text{for } j \neq \operatorname{argmax}_i x_i \end{cases}$$

Choice of loss function

MSE

$$\mathcal{L} = \frac{1}{n} \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2$$

L2

$$\mathcal{L} = \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2$$

Choice of loss function

Binary cross-entropy

$$\mathcal{L} = -\frac{1}{n} \sum_{i=1}^n [y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})]$$

Multiclass:

$$-\frac{1}{n} \sum_n \sum_k y_{nk} \log f_k(\mathbf{x}_n)$$

$$\mathcal{L} = y \log(\sigma(\mathbf{z})) + (1 - y) \log(1 - \sigma(\mathbf{z})),$$

$$\frac{\partial \mathcal{L}}{\partial \theta} = (y - \sigma(\mathbf{z})) \cdot \mathbf{x}$$

Choice of loss function

Binary cross-entropy

$$\mathcal{L} = -\frac{1}{n} \sum_{i=1}^n [y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})]$$

Doesn't have saturation problem

$$\mathcal{L} = y \log(\sigma(\mathbf{z})) + (1 - y) \log(1 - \sigma(\mathbf{z})).$$

$$\frac{\partial \mathcal{L}}{\partial \theta} = (y - \sigma(\mathbf{z})) \cdot \mathbf{x}$$

Deep Networks

Deep networks: informal term for more recent generation of neural nets, with features such as:

- more hidden layers
- built from more heterogenous units
 - sigmoid, rectilinear, max pooling, LSTM, ...
- shared weights across units (convolutional)
- with application-specific network architecture
 - time series, computer vision, speech recognition, ...
 - recurrent networks, max-pooled convolutional layers with local receptive fields...
- bi-directional units
- pretrained on unlabelled data (auto-encoders)
- ...

Impact of Deep Learning

- Speech Recognition



- Computer Vision



- Recommender Systems



- Language Understanding

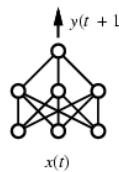
- Drug Discovery and Medical Image Analysis



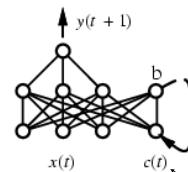
[Courtesy of R. Salakhutdinov]

Recurrent Networks: Time Series

- Suppose we want to predict next state of world
 - and it depends on history of unknown length
 - e.g., robot with forward-facing sensors trying to predict next sensor reading as it moves and turns
- Idea: use hidden layer in network to capture/remember state history



(a) Feedforward network



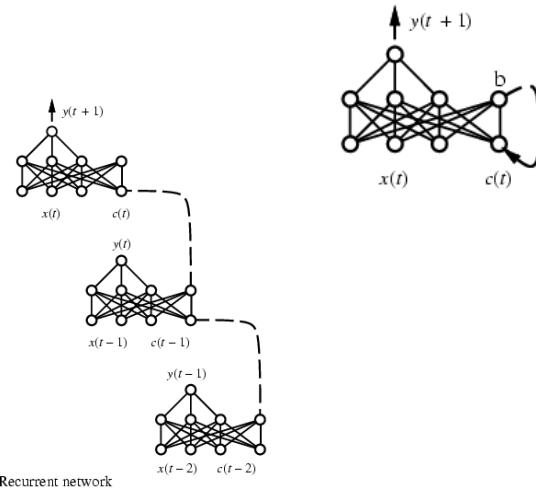
(b) Recurrent network

context/history

Training Networks on Time Series

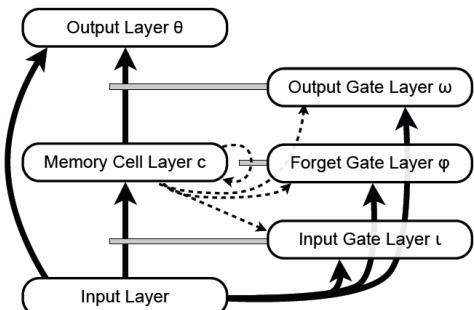
- Suppose we want to predict next state of world
 - and it depends on history of unknown length
 - e.g., robot with forward-facing sensors trying to predict next sensor reading as it moves and turns
 - e.g., anticipate the next word in the sentence

Recurrent Networks on Time Series



Long-Short Term Memory (LSTM) Units

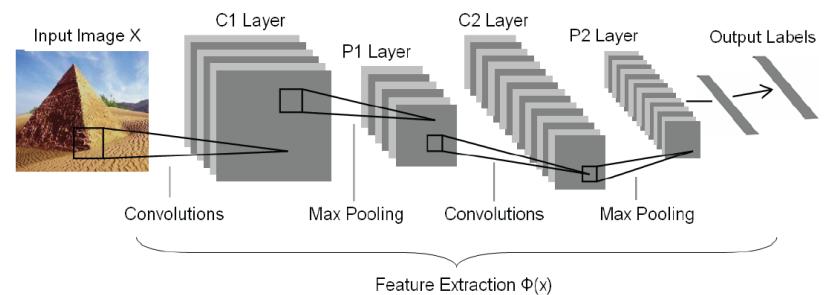
- Threshold unit/subnetwork with memory
 - still trainable with gradient descent



LSTM-g unit from [Monner & Regia, 2013]

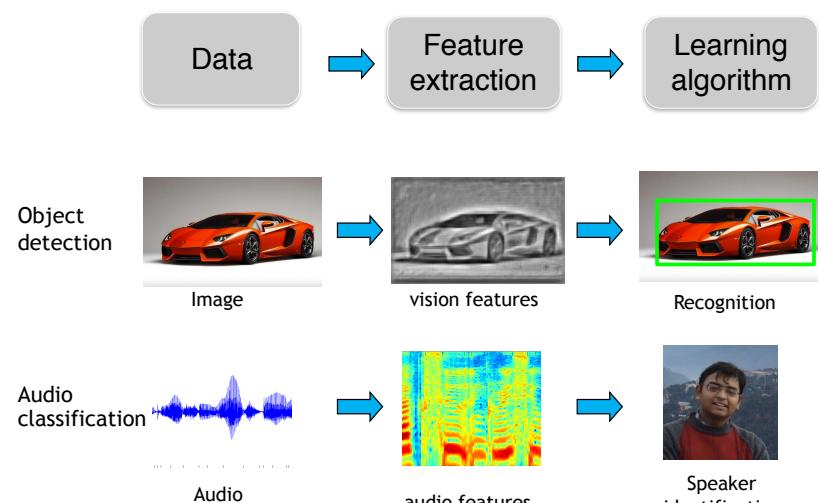
Convolutional Neural Nets for Image Recognition

[Le Cun, 1992]



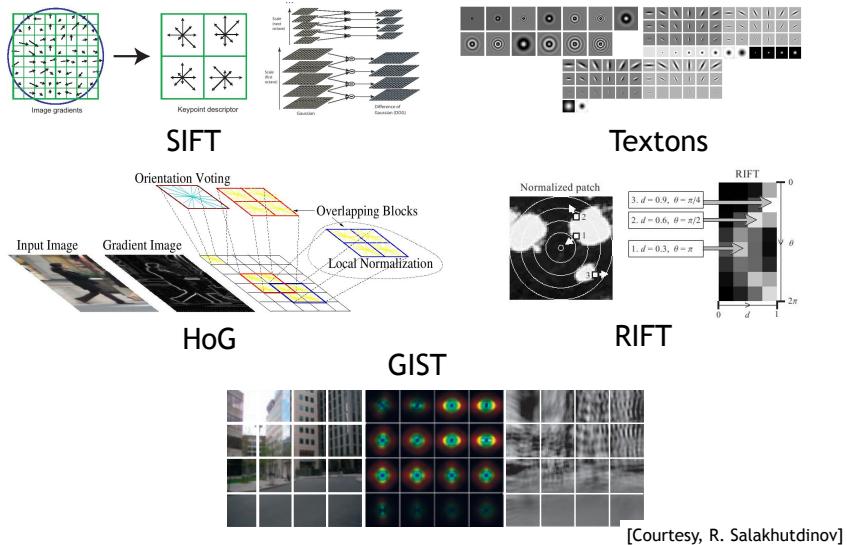
- specialized architecture: mix different types of units, not completely connected, motivated by primate visual cortex
- many shared parameters, stochastic gradient training
- very successful! now many specialized architectures for vision, speech, translation, ...

Feature Representations: Traditionally

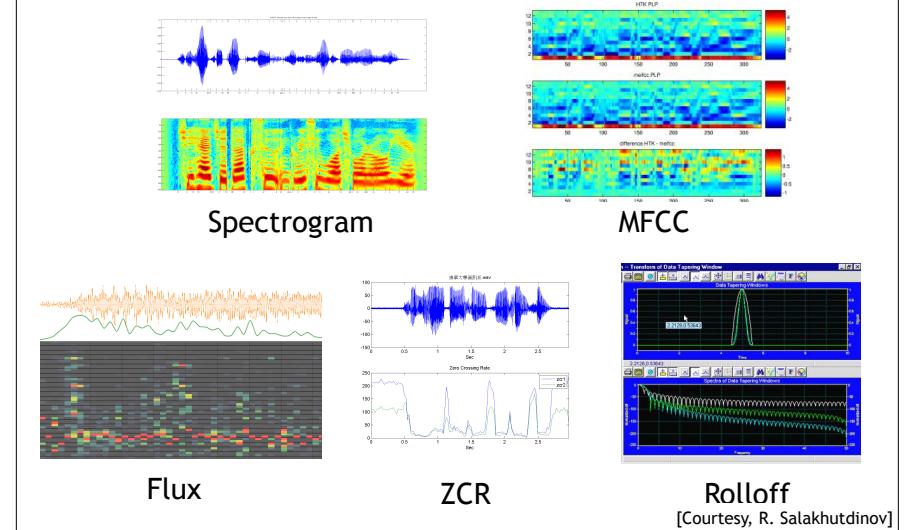


[Courtesy of R. Salakhutdinov]

Computer Vision Features



Audio Features



Audio Features

